

PCA of  $\Sigma$ : performing the computations we obtain:

$$\begin{cases} \lambda_1 = 100.16 & , \quad e_1 = (0.04, 0.999)' \\ \lambda_2 = 0.84 & , \quad e_2 = (0.999, -0.04)' \end{cases}$$

Then we have:

$$\begin{cases} Y_1 = 0.04 X_1 + 0.999 X_2. \rightarrow \text{basically } X_2. \\ Y_2 = 0.999 X_1 + 0.04 X_2. \rightarrow \text{basically } X_1. \end{cases}$$

An example is the dataset of the competitions in athletics of various nations. Using PCA on that dataset, we would find that the PCs will be the marathon.

PCA of  $\mathcal{S}$ :

$$\begin{cases} \lambda_1 = 4.4 & , \quad e_1 = \left( \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right)' \\ \lambda_2 = 0.6 & , \quad e_2 = \left( \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right)' \end{cases}$$

$$\begin{cases} Y_1 = \frac{1}{\sqrt{2}} Z_1 + \frac{1}{\sqrt{2}} Z_2 = \frac{1}{\sqrt{2}} \left( \frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}} + \frac{X_2 - \mu_2}{\sqrt{\sigma_{22}}} \right) \\ Y_2 = \frac{1}{\sqrt{2}} Z_1 - \frac{1}{\sqrt{2}} Z_2 = \frac{1}{\sqrt{2}} \left( \frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}} - \frac{X_2 - \mu_2}{\sqrt{\sigma_{22}}} \right). \end{cases}$$

The linear combinations forming  $Y_1, Y_2$  are very different from the ones of before!

THM: if data is very influenced by units of measure etc., then use the standardized data  $Z$  and apply PCA on it.

Summarizing Sample Variation by Principal Components:

Usually,  $\mu$  and  $\Sigma$  are unknown. But we have data!

$$X = \begin{bmatrix} x_1' \\ \vdots \\ x_n' \end{bmatrix} \quad \text{data, } x_i \text{ realization of } X_i, \\ x_1, \dots, x_n \text{ iid } \sim \mathcal{X}.$$

Hence we will use data to estimate  $\mu$  and  $\Sigma$ .

In particular,  $\Sigma$  is estimated by  $S$  and  $\mu$  by  $\bar{x}$ .

$\Rightarrow$  We perform PCA on  $S$ :

$$S = \sum_{i=1}^p \lambda_i e_i e_i' \Rightarrow \text{PC } Y_i: \text{projection on } e_i.$$

$$x_i \xrightarrow{\text{PCA}} y_i = \begin{pmatrix} (e_1' x_i) \\ \vdots \\ (e_p' x_i) \end{pmatrix} \quad \text{scores of } x_i \text{ (i.e. projection of } x_i \text{ in the space of eigenvectors } P)$$

Hence:

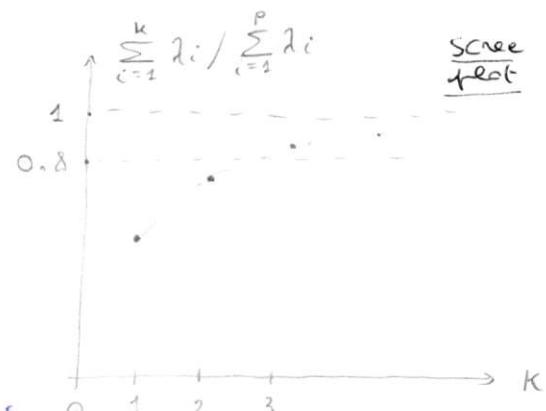
$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \xrightarrow{\text{PCA}} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \dots & y_{np} \end{bmatrix}.$$

Number of principal components:

Suppose that you set a threshold:

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^p \lambda_i} \geq \text{threshold} \quad (\text{e.g. } 0.8).$$

In this way, we can reduce the analysis to the first  $k$  principal components.



Observation 1: PCA for categorical variables

is called correspondence analysis.

It is performed on the table of joint frequencies (contingency tables).

Observation 2: look also at the smallest  $\lambda_i$ .

If  $\lambda_p \approx 0$ , it means that there is a linear relationship between  $x_1, \dots, x_k$ .

We now look at a more "mathematical fashion" for the derivation of PCA.

Derivation of PCA: Optimal orthonormal basis:

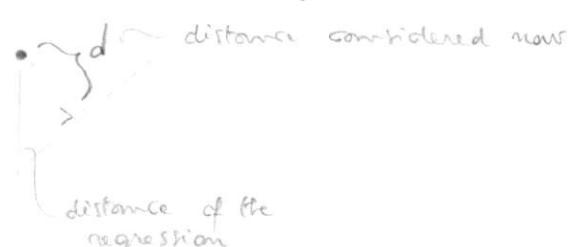
$x_1, \dots, x_n \in \mathbb{R}^p$ .

Problem: find a linear space  $L$  of dimension  $K \leq p$  s.t.  $L$  is "closest" to the data.

Example:  $p=2, K=1$ .



What do we mean by closest?



$\Rightarrow$  that line is not the regression line!

That line is a geometrical property: the minimization of the orthogonal distance doesn't depend on the axis we've chosen, while regression depends on the chosen basis.

Let  $\{\eta_1, \dots, \eta_K\}$  orthonormal basis spanning  $\mathcal{L}$ , so that:

$$\mathcal{L} = \text{span} \{ \eta_1, \dots, \eta_K \}.$$

Problem: find  $\eta_1, \dots, \eta_K$  o.e. s.t.

$$\sum_{i=1}^n \|(\underline{x}_i - \bar{x}) - \sum_{j=1}^K \eta_j \eta_j' (\underline{x}_i - \bar{x})\|^2 \text{ is } \underline{\min}.$$

↓ of before, call it (a).

Let  $\underline{v}_i := \underline{x}_i - \bar{x}$ .

$$\begin{aligned} (a) &= \|\underline{v}_i - \sum_{j=1}^K \eta_j \eta_j' \underline{v}_i\|^2 = \\ &= \left( \underline{v}_i - \sum_{j=1}^K \eta_j \eta_j' \underline{v}_i \right)' \left( \underline{v}_i - \sum_{j=1}^K \eta_j \eta_j' \underline{v}_i \right) = \\ &= \underline{v}_i' \underline{v}_i - 2 \sum_{j=1}^K \underline{v}_i' \eta_j \eta_j' \underline{v}_i + \left( \sum_{j=1}^K \eta_j \eta_j' \underline{v}_i \right)' \left( \sum_{j=1}^K \eta_j \eta_j' \underline{v}_i \right) = \\ &= \underline{v}_i' \underline{v}_i - 2 \sum_{j=1}^K (\eta_j' \underline{v}_i)^2 + \sum_{j=1}^K (\eta_j' \underline{v}_i)^2 = \\ &= \underline{v}_i' \underline{v}_i - \sum_{j=1}^K (\eta_j' \underline{v}_i)^2. \end{aligned}$$

Then,  $\sum_{i=1}^n \left( \underline{v}_i' \underline{v}_i - \sum_{j=1}^K (\eta_j' \underline{v}_i)^2 \right)$  is min  $\iff$

$$\sum_{i=1}^n \sum_{j=1}^K (\eta_j' \underline{v}_i)^2 \text{ is } \underline{\max}.$$

$$\begin{aligned} &= \sum_{j=1}^K \left( \sum_{i=1}^n \eta_j' \underline{v}_i \underline{v}_i' \eta_j \right) = \\ &= \sum_{j=1}^K \eta_j' \left( \underbrace{\sum_{i=1}^n (\underline{x}_i - \bar{x})(\underline{x}_i - \bar{x})'}_{=(n-1)S} \right) \eta_j = \\ &= (n-1) \sum_{j=1}^K \eta_j' S \eta_j. \end{aligned}$$

K=1:  $\max_{\eta: \|\eta\|=1} \eta' S \eta = \lambda_1$ ,  $\arg \max \dots = \underline{e}_1$ .

By induction,  $K=2 = \lambda_2$  and so on.

$$\Rightarrow \boxed{\eta_1 = \underline{e}_1, \dots, \eta_K = \underline{e}_K} \text{ s.t. } S = \sum_i \lambda_i \underline{e}_i \underline{e}_i'.$$

$\rightarrow \{\eta_1, \dots, \eta_K\}$  are the first K eigenvectors of  $S$ : we have found an alternative derivation of PCA!

(end)

### Observation:

$$\sum_{i=1}^n \text{det} \left( \underline{x}_i - \bar{\underline{x}}, \sum_{j=1}^k e_j e_j' (\underline{x}_i - \bar{\underline{x}}) \right) =$$

$$= \underbrace{\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})' (\underline{x}_i - \bar{\underline{x}})}_{\in \mathbb{R}} - (n-1) \sum_{j=1}^k e_j' S e_j = [...].$$

Now,

$$\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})' (\underline{x}_i - \bar{\underline{x}}) = \text{tr} \left( \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})' (\underline{x}_i - \bar{\underline{x}}) \right) =$$

$$= - \sum_{i=1}^n \text{tr} \left( (\underline{x}_i - \bar{\underline{x}})' (\underline{x}_i - \bar{\underline{x}}) \right) \stackrel{(*)}{=} .$$

$$= \text{tr} \left( \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}}) (\underline{x}_i - \bar{\underline{x}})' \right) =$$

$$= \text{tr} ((n-1) S) = (n-1) \sum_{i=1}^p \lambda_i. \quad \text{Moreover, } (n-1) \sum_{j=1}^k e_j' S e_j =$$

$$= (n-1) \sum_{i=1}^k \lambda_i.$$

$$\Rightarrow [...] = (n-1) \sum_{i=1}^p \lambda_i - (n-1) \sum_{j=1}^k \lambda_j$$

$$= (n-1) \sum_{j=k+1}^p \lambda_j. \quad \text{error of approximation (the dimension of the orth. basis is } k, \text{ not } p)$$

The PCA was introduced in this way by Pearson in 1900.

### Extensions of PCA:

- ICA: Independent Component Analysis.
- Non-linear dimensional reduction. Very open field, we are entering the machine-learning territory.

# Multivariate Gaussian Distribution (Ch.4, JB):

24/3/20

$\underline{X} \in \mathbb{R}^P$  random vector.

$\underline{\gamma}_X : \mathcal{B}$  (Borel sets of  $\mathbb{R}^P$ )  $\rightarrow [0, 1]$

$$\underline{\gamma}_X(B) = P(\underline{X} \in B) \quad \forall B \subseteq \mathcal{B}.$$

In most cases, we can write  $\underline{\gamma}_X$  as:

$$\underline{\gamma}_X(B) = \int_B f(\underline{x}) d\underline{x}, \quad f: \text{density}.$$

Properties of  $f$ :

$$1) f \geq 0.$$

$$2) \int_{\mathbb{R}^P} f(\underline{x}) d\underline{x} = 1.$$

Definition:  $\underline{X} \in \mathbb{R}^P$  random vector with density  $f$ .

Let  $\mu \in \mathbb{R}^P$ ,  $\Sigma$   $p \times p$  positive definite.

Then  $\underline{X}$  is said to have Gaussian (or Normal) distribution with parameters  $(\mu, \Sigma)$ , namely  $\underline{X} \sim N_p(\mu, \Sigma)$ , if:

$$f(\underline{x}) = \frac{1}{\sqrt{(2\pi)^P \det \Sigma}} \exp\left(-\frac{1}{2} (\underline{x} - \mu)^T \Sigma^{-1} (\underline{x} - \mu)\right).$$

$$\propto \exp\left(-\frac{1}{2} d_{\Sigma^{-1}}^2(\underline{x}, \mu)\right).$$

Mahalanobis distance

probability  
density function  
of a multi-  
variate Gaussian

Contour plots for  $f$  Gaussian:

$$\{\underline{x} \in \mathbb{R}^P : f(\underline{x}) = \text{const.}\} =$$

$$= \{\underline{x} \in \mathbb{R}^P : (\underline{x} - \mu)^T \Sigma^{-1} (\underline{x} - \mu) = \text{const}^2\} =$$

$$= \{\underline{x} \in \mathbb{R}^P : d_{\Sigma^{-1}}^2(\underline{x}, \mu) = \text{const}^2\}.$$

Writing the spectral decomposition for  $\Sigma$ :

$$\Sigma = \sum_{i=1}^P \lambda_i e_i e_i^T \Rightarrow \Sigma^{-1} = \sum_{i=1}^P \frac{1}{\lambda_i} e_i e_i^T.$$

Then, the axes of the ellipses will be:  $e_1, \dots, e_P$ ,  
while the length of the axes will be:  $\propto \sqrt{\lambda_1}, \dots, \sqrt{\lambda_p}$ .



X1

Proposition 1: Let  $\underline{X} \sim N_p(\mu, \Sigma)$ . Then:

$$E[\underline{X}] = \mu, \quad \text{cov}(\underline{X}) = \Sigma.$$

Theorem 2:  $\underline{X} \sim N_p(\mu, \Sigma) \Leftrightarrow \underline{a}' \underline{X} \sim N_1(a'\mu, a'\Sigma a) \quad \forall a \in \mathbb{R}^p$ .

This theorem gives us a characterization of the Gaussian distribution.

Proof: use the characteristic function of the Gaussian distribution. (end)

The theorem has a strong impact, both practically and theoretically.

Corollary:

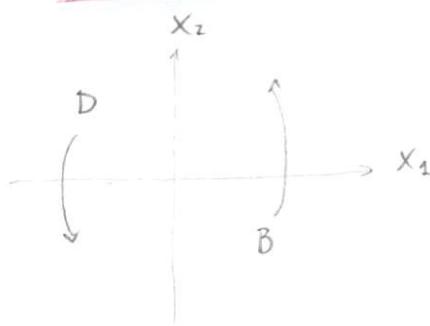
$$\underline{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix}, \quad \underline{X} \sim N_p(\mu, \Sigma) \Rightarrow X_i \sim N_1(\mu_i, \sigma_{ii}).$$
$$\Sigma = [\sigma_{ij}]$$

Hence, Gaussianity of a vector implies Gaussianity of its components. The vice-versa is generally wrong.

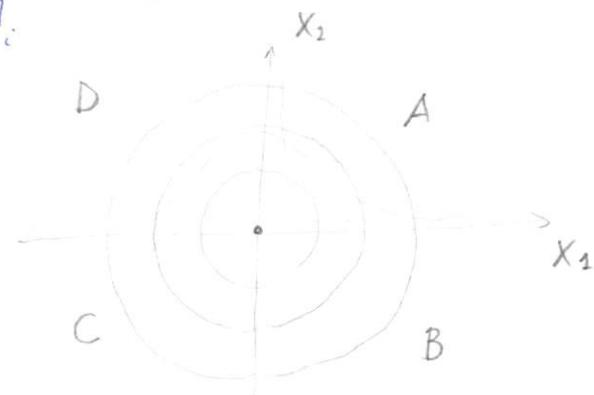
Indeed, a component is a linear combination of  $\underline{X}$ :

$$X_i = u_i' \underline{X}, \quad u_i = (0 \dots 0 \underset{i}{\underset{\downarrow}{1}} 0 \dots 0)'.$$

Example: ( $\times$ ).  $\underline{X} \sim N_2(0, I)$ .



We get the following distribution:



This is not anymore Gaussian!

Gaussianity of components is necessary, but not sufficient!

Proposition 2:  $\underline{X} \sim N_p(\underline{\mu}, \Sigma)$ ,  $A$   $q \times p$  matrix. Then:

$$A\underline{X} \in \mathbb{R}^q \sim N_q(A\underline{\mu}, A\Sigma A').$$

The important fact of this proposition is that Generativity is preserved.

Proof:

Let  $\underline{a} \in \mathbb{R}^q$ . We need to prove that  $\underline{a}'(A\underline{X})$  has G. distn.

$$\underline{a}'(A\underline{X}) = (\underline{a}'A)\underline{X} = (A'\underline{a})'\underline{X}.$$

But  $A'\underline{a} \in \mathbb{R}^p$  <sup>(trans.)</sup>  $\Rightarrow (A'\underline{a})'\underline{X} \sim N_1((A'\underline{a})'\underline{\mu}), (A'\underline{a})'\Sigma(A'\underline{a}) \sim N_1(\underline{a}'(A\underline{\mu}), \underline{a}'A\Sigma A'\underline{a}) \forall \underline{a} \in \mathbb{R}^q$ .

<sup>(trans)</sup>  
 $\Rightarrow A\underline{X} \sim N_q(A\underline{\mu}, A\Sigma A')$ . (end)

Proposition 3:  $\underline{X} \sim N_p(\underline{\mu}, \Sigma)$ ,  $d \in \mathbb{R}^p$ .

$$\Rightarrow \underline{X} + d \sim N_p(\underline{\mu} + d, \Sigma).$$

Proof: as exercise. Prove it as before. (end)

An even more interesting thing is the following. Let:

$$\underline{X} = \begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix}, \quad \underline{X}_1 \in \mathbb{R}^q, \quad \underline{X}_2 \in \mathbb{R}^{p-q}, \quad q < p.$$

$$\text{Let } \underline{\mu} = \begin{pmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{pmatrix}, \quad \underline{\mu}_1 \in \mathbb{R}^q, \quad \underline{\mu}_2 \in \mathbb{R}^{p-q},$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} (q \times q) & \Sigma_{12} (q \times (p-q)) \\ \Sigma_{21} ((p-q) \times q) & \Sigma_{22} ((p-q) \times (p-q)) \end{bmatrix} \quad (\text{block notation})$$

Proposition 4: Set  $\underline{X} = \begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix} \sim N_p(\underline{\mu}, \Sigma)$ .

Then:  $\underline{X}_1 \sim N_q(\underline{\mu}_1, \Sigma_{11})$ .

Proof: Again, using linear transformations.

$$\text{Let } A = [I_{(q \times q)} \quad 0_{q \times (p-q)}].$$

<sup>thus</sup>

$$A\underline{X} \sim N_q(A\underline{\mu}, A\Sigma A').$$

$$\text{But } A\underline{X} = \underline{X}_1, \quad A\underline{\mu} = \underline{\mu}_1, \quad A\Sigma A' = \Sigma_{11}. \quad (\text{end})$$

Theorem 2: let  $\underline{X} = \begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 \end{pmatrix} \sim N_p \left( \begin{pmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{pmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$ . Then:

$$\underline{x}_1 \perp\!\!\!\perp \underline{x}_2 \iff \Sigma_{12} = 0 \quad (\text{by symmetry, } \Sigma_{21} = 0).$$

This theorem is telling us that, in the Gaussian world, covariance equal to zero is equivalent to independence.

Proof: we can split the density of  $f_{\underline{X}}$  into:

$$f_{\underline{X}} = f_{\underline{x}_1} \cdot f_{\underline{x}_2}. \quad (\text{end})$$

Theorem 3: let  $\underline{X} = \begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 \end{pmatrix} \sim N_p \left( \begin{pmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$ . Then:

$$\underline{x}_1 \mid \underline{x}_2 = \underline{x}_2 \sim N_q \left( \underline{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right) \quad (\text{with } \det \Sigma_{22} \neq 0).$$

Remark: Taking  $\underline{z}_1, \dots, \underline{z}_p \sim \text{iid } N_1(0, 1)$ ,  $\underline{Z} := (\underline{z}_1, \dots, \underline{z}_p)'$ ,

then:  $\underline{X} := \Sigma^{1/2} \underline{Z} + \underline{\mu} \sim N_p(\underline{\mu}, \Sigma)$ .

Proof (of Thm.3):

$$\text{Take: } A = \begin{bmatrix} I_{(q \times q)} & -\Sigma_{12} \Sigma_{22}^{-1} (q \times (p-q)) \\ 0 ((p-q) \times q) & I_{(p-q) \times (p-q)} \end{bmatrix}.$$

$$A(\underline{X} - \underline{\mu}) = A \begin{pmatrix} \underline{x}_1 - \underline{\mu}_1 \\ \underline{x}_2 - \underline{\mu}_2 \end{pmatrix} = \begin{pmatrix} \underline{x}_1 - \underline{\mu}_1 - \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2) \\ \underline{x}_2 - \underline{\mu}_2 \end{pmatrix}$$

$$\sim N_p \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, A \Sigma A' \right) \sim N_p \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} & 0 \\ 0 & \Sigma_{22} \end{pmatrix} \right).$$

$$\underline{x}_1 - \underline{\mu}_1 - \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2) \perp\!\!\!\perp \underline{x}_2 - \underline{\mu}_2.$$

$$\text{Moreover, } \underbrace{\underline{x}_1 - \underline{\mu}_1 - \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2)}_{W} \sim N_q(0, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}).$$

$$\text{Now, } W \mid \underline{x}_2 = \underline{x}_2 \sim N_q(0, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}) \sim \text{because of independence}$$

$$\Rightarrow \underline{x}_1 - \underline{\mu}_1 - \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2) \sim N_q(0, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}).$$

$$\text{By translation, } \underline{x}_1 \mid \underline{x}_2 = \underline{x}_2 \sim N_q(\underline{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}). \quad (\text{end})$$

Remark:  $\text{Cov}(X_1 | X_2 = x_2) = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$  }  
does not depend on  $x_2$ . } partial covariances

Example:  $p = 2$ .

$$X = \begin{pmatrix} X \\ Y \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \begin{pmatrix} \sigma_{XX} & \sigma_{XY} \\ \sigma_{YX} & \sigma_{YY} \end{pmatrix} \right).$$

For example, height and weight ( $Y$ ).

$$Y \sim N(\mu_Y, \sigma_{YY}) \Rightarrow P[Y \in [\mu_Y \pm 2\sqrt{\sigma_{YY}}]] = 0.95.$$

$$Y | X=x \sim N_1 \left( \mu_Y + \sigma_{XY} \sigma_{XX}^{-1} (x - \mu_X), \sigma_{YY} - \sigma_{XY} \sigma_{XX}^{-1} \sigma_{YX} \right) = \\ \sim N_1 \left( \mu_Y + \frac{\sigma_{XY}}{\sigma_{XX}} (x - \mu_X), \sigma_{YY} - \frac{\sigma_{XY}^2}{\sigma_{XX}} \right).$$

$$\text{Recall that } \rho_{XY} = \frac{\sigma_{XY}}{\sqrt{\sigma_{XX} \sigma_{YY}}} \in [-1, 1].$$

$$\sim N_1 \left( \mu_Y + \frac{\sigma_{XY}}{\sigma_{XX}} (x - \mu_X), \sigma_{YY} (1 - \rho^2) \right).$$

if there is some correlation between  $X$  and  $Y$ , the prediction on  $Y$  will be more accurate (less variable).

$$y = E[Y | X=x] = \mu_Y + \frac{\sigma_{XY}}{\sigma_{XX}} (x - \mu_X).$$

↳ linear regression function

Rank: In the Gaussian world, the regression function is exactly linear.

This can be rewritten as:

$$\frac{y - \mu_Y}{\sqrt{\sigma_{YY}}} = \frac{\sigma_{XY}}{\sqrt{\sigma_{YY} \cdot \sigma_{XX}}} \cdot \frac{x - \mu_X}{\sqrt{\sigma_{XX}}} .$$

$\rho_{XY}$

$$\Rightarrow \frac{y - \mu_Y}{\sqrt{\sigma_{YY}}} = \rho_{XY} \frac{x - \mu_X}{\sqrt{\sigma_{XX}}} .$$

standardized data

We see that an exceptionality on  $X$  (like  $x = \mu_X + 2\sqrt{\sigma_{XX}}$ ) is not an exceptionality on  $Y$  ( $y = \mu_Y + 2\rho\sqrt{\sigma_{YY}}$ ).

History: Sir Galton was regressing the height of fathers vs. sons.  
 $\Rightarrow$  "regressing towards mediocrity".

Regression effect

$$|\rho| < 1.$$

The regression effect causes regression fallacy.

Example: notes on calculus I and calculus II.

Note on CI : 30 → prediction: CII 28.

Note on CI : 18 → prediction: CII 20.

The fallacy is to say: the smart guy will study less in the next exam, while the other guy was disappointed and will study more because ... ⇒ wrong reasoning.

Suggestion: Daniel Kahneman, "Thinking Fast & Slow".

Homework: read the two articles about regression fallacy.

26/3/20

Let  $\underline{X}$  r.v.  $\sim N_p(\mu, \Sigma)$ .

Proposition: If  $\det(\Sigma) > 0$ , then:

$$(\underline{X} - \mu)' \Sigma^{-1} (\underline{X} - \mu) \sim \chi^2(p).$$

Proof:

Remember: if  $z_1, \dots, z_p$  are i.i.d  $\sim N(0, 1)$ , then

$$\sum_{i=1}^p z_i^2 \sim \chi^2(p).$$

Consider  $\Sigma = \sum_i \lambda_i e_i e_i'$  spectral decomposition,

$$P = [e_1 \dots e_p], \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p) \Rightarrow \Sigma = P \Lambda P'.$$

Consider:  $\Lambda^{-\frac{1}{2}} P' (\underline{X} - \mu) \sim N_p(0, \Lambda^{-\frac{1}{2}} P \Sigma P' \Lambda^{-\frac{1}{2}}).$

$$\text{Now, } \Lambda^{-\frac{1}{2}} P \Sigma P' \Lambda^{-\frac{1}{2}} =$$

$$= \Lambda^{-\frac{1}{2}} P' P \Lambda P' P \Lambda^{-\frac{1}{2}} = \Lambda^{-\frac{1}{2}} \Lambda \Lambda^{-\frac{1}{2}} = I.$$

$$\Rightarrow z = \Lambda^{-\frac{1}{2}} P' (\underline{X} - \mu) \sim N_p(0, I).$$

$$z = \begin{pmatrix} z_1 \\ \vdots \\ z_p \end{pmatrix} \Rightarrow z_1 \dots z_p \text{ iid } \sim N(0, 1).$$

I can write:

$$(\underline{X} - \mu)' \Sigma^{-1} (\underline{X} - \mu) = (\underline{X} - \mu)' P \Lambda^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} P' (\underline{X} - \mu) = z' z =$$

$$= \sum_{i=1}^p z_i^2 \sim \chi^2(p).$$

(end)

Observation: if  $\det(\Sigma) = 0$ , let  $K = \text{rank}(\Sigma)$ .

$$\Sigma = \sum_{i=1}^p \lambda_i e_i e_i^\top, \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_K > 0 = \lambda_{K+1} = \dots = \lambda_p.$$

We can consider a "generalized inverse" for  $\Sigma$ :

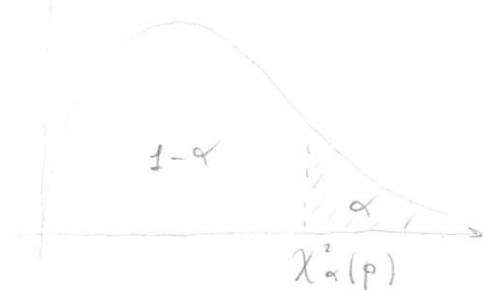
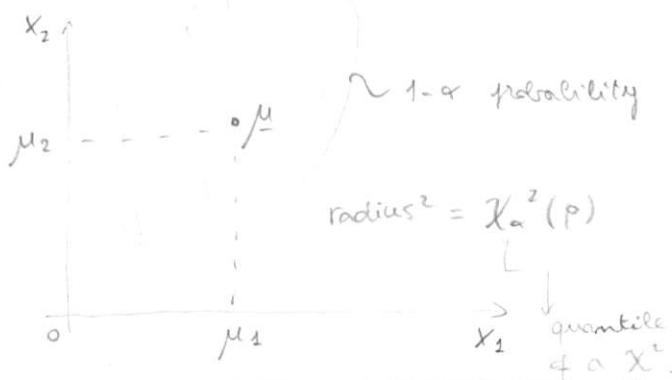
$$\Sigma^- := \sum_{i=1}^K \frac{1}{\lambda_i} e_i e_i^\top.$$

generalized inverse for non-full-rank  $\Sigma$ .

Prove that:  $(\underline{X} - \mu)^\top \Sigma^- (\underline{X} - \mu) \sim \chi^2(K)$ .

Corollary: Let  $\alpha \in (0, 1)$ . Then ( $\text{Det}(\Sigma) > 0$ ):

$$\Pr[(\underline{X} - \mu)^\top \Sigma^{-1} (\underline{X} - \mu) \leq \chi_{\alpha}^2(p)] = 1 - \alpha. \quad \begin{matrix} \text{confidence region} \\ (\text{level } \alpha) \text{ for } \underline{X} \sim N_p(\mu, \Sigma) \end{matrix}$$



Remark: in R,  $\chi_{\alpha}^2(p) = \text{qchisq}(1 - \alpha, p)$ .

Estimators of  $\mu$  and  $\Sigma$  for  $N_p(\mu, \Sigma)$ :

$\underline{X} = \begin{bmatrix} \underline{x}_1^\top \\ \vdots \\ \underline{x}_n^\top \end{bmatrix}$  data,  $\underline{x}_i^\top \xleftarrow{\text{obs.}} X_i$  random vector.  
Hp:  $X_1, \dots, X_n$  iid  $\sim N_p(\mu, \Sigma)$ ,  
 $\mu, \Sigma$  unknown.

Obvious choices:

$$\begin{aligned} \cdot \bar{X} &= \frac{1}{n} \sum_{i=1}^n \underline{x}_i \\ \cdot S &= \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{X})(\underline{x}_i - \bar{X})^\top \end{aligned} \quad \left. \right\} \text{unbiased estimators for } \mu, \Sigma.$$

MLE estimators for  $\mu$  and  $\Sigma$ :

Suppose we want to compute the probability:

$$\text{joint density } [X_1 = \underline{x}_1, \dots, X_n = \underline{x}_n] = \prod_{i=1}^n \frac{1}{\sqrt{(2\pi)^p \text{Det}(\Sigma)}} \exp\left[-\frac{1}{2}(\underline{x}_i - \mu)^\top \Sigma^{-1} (\underline{x}_i - \mu)\right].$$

Let's define now a function:

$$L: (\mu, \Sigma) \rightarrow \text{joint density } [X_1 = \underline{x}_1, \dots, X_n = \underline{x}_n]$$

given  $\underline{X}_1 = \underline{x}_1, \dots, \underline{X}_n = \underline{x}_n$ .

$L$  is called Likelihood function.

Theorem: Maximum Likelihood Estimators for  $\mu, \Sigma$ :

$$\arg \max_{(\mu, \Sigma)} L(\mu, \Sigma | \underline{X}_1 = \underline{x}_1, \dots, \underline{X}_n = \underline{x}_n) = (\hat{\mu}, \hat{\Sigma}).$$

$(\mu, \Sigma) : \mu \in \mathbb{R}^p,$   
 $\Sigma p \times p \text{ pos. def.}$

where  $L(\mu, \Sigma | \underline{X}_1 = \underline{x}_1, \dots, \underline{X}_n = \underline{x}_n) = \prod_{i=1}^n \frac{1}{\sqrt{(2\pi)^p \det \Sigma}} \exp \left[ -\frac{1}{2} (\underline{x}_i - \mu)' \Sigma^{-1} (\underline{x}_i - \mu) \right],$

and  $\begin{cases} \hat{\mu} = \bar{\underline{x}}, \\ \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})' (\underline{x}_i - \bar{\underline{x}}) = \frac{n-1}{n} S. \end{cases}$  MLE's for  $(\mu, \Sigma)$

General Properties of MLE's:

- Let  $\theta \in \mathbb{R}^K$  parameter (ex.  $\theta = (\mu, \Sigma)$ ),  $\hat{\theta} = \hat{\theta}(\text{data})$  MLE estimator of  $\theta$ . Let  $h: \mathbb{R}^K \rightarrow \mathbb{R}^j$ .

What is the MLE of  $h(\theta)$ ? This is written  $\hat{h}(\hat{\theta})$ .

The invariance property of MLE holds:

$$\hat{h}(\hat{\theta}) = h(\hat{\theta}).$$

Example:  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})' (\underline{x}_i - \bar{\underline{x}})$ , MLE for  $\Sigma$   
if  $\underline{X}_1, \dots, \underline{X}_n$  iid  $\sim N_p(\mu, \Sigma)$ .

Which is a good estimator of  $\lambda_1$ , where  $\Sigma = \sum_i \lambda_i \underline{e}_i \underline{e}_i'$ ?  
By the property of before,  $\hat{\lambda}_1$  is s.t.  $\hat{\Sigma} = \sum_i \hat{\lambda}_i \underline{e}_i \underline{e}_i'$ .

Sampling distribution of  $\bar{\underline{x}}$  and  $S$  (or  $\hat{\Sigma}$ ):

Assume  $\underline{X}_1, \dots, \underline{X}_n$  iid  $\sim N_p(\mu, \Sigma)$ .

Proposition:  $\bar{\underline{x}} \sim N_p(\mu, \frac{\Sigma}{n})$ .

distribution of  $\bar{\underline{x}} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i$

Proof:  $\tilde{\underline{x}} = \begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_n \end{pmatrix} \in \mathbb{R}^{np}$   $n \times 1$  matrix.

Then (iid)  $\tilde{\underline{x}} \sim N_{np} \left( \begin{pmatrix} \mu \\ \mu \\ \vdots \\ \mu \end{pmatrix}, \begin{pmatrix} \Sigma & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \Sigma \end{pmatrix} \right).$

Consider  $A = [\underbrace{I(p \times p)}_n \ I \ \dots \ I]^T$ .  $p \times np$ .

$$\Rightarrow \frac{1}{n} A \tilde{X} = \bar{X} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i.$$

Indeed,  $A = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$

$$\Rightarrow \frac{1}{n} A \tilde{X} = \frac{1}{n} \begin{pmatrix} \sum_{i=1}^n x_{i1} \\ \vdots \\ \sum_{i=1}^n x_{ip} \end{pmatrix} = \bar{X} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i.$$

Hence:  $\bar{X} = \frac{1}{n} A \tilde{X} \sim N_p \left( \frac{1}{n} A \begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix}, \frac{1}{n^2} A \begin{pmatrix} \Sigma & & 0 \\ & \ddots & \\ 0 & & \Sigma \end{pmatrix} A^T \right)$   
 $\sim N_p \left( \mu, \frac{1}{n^2} n \Sigma \right) \sim N_p \left( \mu, \frac{1}{n} \Sigma \right).$

(end)

► What about the distribution of  $S$ ? We are talking about a random positive-semidefinite matrix.

Definition: Let  $\underline{z}_1, \dots, \underline{z}_m$  iid  $\sim N_p(\underline{\sigma}, \Sigma)$ ,  $\det(\Sigma) > 0$ .

Then  $\sum_{i=1}^m \underline{z}_i \underline{z}_i^T \sim \text{Wishart}(\Sigma, m)$

Wishart Distribution  
 $m$ : degrees of freedom

Properties of Wishart distribution:

1)  $A_1 \sim \text{Wish}(\Sigma, m_1)$ ,  $A_2 \sim \text{Wish}(\Sigma, m_2)$ ,  $A_1 \perp\!\!\!\perp A_2$ .  
 (stoch. independent)

$$\Rightarrow A_1 + A_2 \sim \text{Wish}(\Sigma, m_1 + m_2). \quad \text{sums of Wisharts}$$

Proof:  $\begin{cases} A_1 = \sum_{i=1}^{m_1} \underline{z}_i \underline{z}_i^T, & \underline{z}_1, \dots, \underline{z}_{m_1} \text{ iid } \sim N_p(\underline{\sigma}, \Sigma) \\ A_2 = \sum_{i=1}^{m_2} \tilde{\underline{z}}_i \tilde{\underline{z}}_i^T, & \tilde{\underline{z}}_1, \dots, \tilde{\underline{z}}_{m_2} \text{ iid } \sim N_p(\underline{\sigma}, \Sigma). \end{cases}$

Renaming the rvs:

$$\begin{array}{cccccc} \underline{z}_1, \dots, & \underline{z}_{m_1}, & \tilde{\underline{z}}_1, \dots, & \tilde{\underline{z}}_{m_2} \\ \parallel & \parallel & \parallel & \parallel \\ W_1 & W_2 & W_{m_1} & W_{m_1+1} & \dots & W_{m_1+m_2} \end{array}$$

$$W_1, \dots, W_{m_1+m_2} \text{ iid } \sim N_p(\underline{\sigma}, \Sigma).$$

$$\Rightarrow A_1 + A_2 = \sum_{i=1}^{m_1} \underline{z}_i \underline{z}_i^T + \sum_{i=1}^{m_2} \tilde{\underline{z}}_i \tilde{\underline{z}}_i^T = \sum_{i=1}^{m_1+m_2} W_i W_i^T \sim \text{Wish}(\Sigma, m_1 + m_2)$$

(end)

2) Let  $C \in \mathbb{R}^{p \times p}$  constant,  $A \sim \text{Wish}(\Sigma, m)$ .

$$\Rightarrow C A C' \sim \text{Wish}(C \Sigma C', m).$$

Proof:  $A = \sum_{i=1}^m \underline{z}_i \underline{z}_i'$ ,  $\underline{z}_i$  iid  $\sim N_p(0, \Sigma)$ .

$$CAC' = \sum_{i=1}^m \underbrace{C \underline{z}_i}_{\underline{w}_i} \underline{z}_i C' , \quad \underline{w}_i \text{ iid } \sim N_p(0, C \Sigma C').$$

$\underline{w}_i \sim \text{Wish}(C \Sigma C', m)$  by definition.  
(end)

3)  $\sigma^2 > 0$ ,  $A \sim \text{Wish}(\Sigma, m)$ .

$$\Rightarrow \sigma^2 A \sim \text{Wish}(\sigma^2 \Sigma, m).$$

Proof:  $\sigma^2 A = \sum_{i=1}^m \underbrace{\sigma \underline{z}_i \underline{z}_i'}_{\underline{w}_i} \sigma$   
 $\underline{w}_i \rightarrow \underline{w}_i \text{ iid } \sim N_p(0, \sigma^2 \Sigma).$

$\sim \text{Wish}(\sigma^2 \Sigma, m)$  by def. (end)

4) Let  $A \sim \text{Wish}(\Sigma, m)$  s.t.  $\Sigma \in \mathbb{R}$  (i.e.  $p=1$ ), ( $\Sigma = \sigma^2$ )

Then:  $A = \sum_{i=1}^m \underline{z}_i \underline{z}_i$ ,  $\underline{z}_1, \dots, \underline{z}_m$  iid  $\sim N_1(0, \Sigma)$ .

$$\Rightarrow \frac{1}{\sigma^2} A = \sum_{i=1}^m \frac{\underline{z}_i}{\sigma} \cdot \frac{\underline{z}_i}{\sigma} \Rightarrow \frac{1}{\sigma^2} A \sim \chi^2(m).$$

$\underline{z}_i \text{ iid } N_1(0, 1)$

In conclusion,

$A \sim \text{Wish}(\Sigma, m)$ ,  $\Sigma \in \mathbb{R}$ .

$$\Rightarrow A \sim \Sigma \cdot \chi^2(m), \quad \text{i.e.} \quad \frac{A}{\Sigma} \sim \chi^2(m). \quad \begin{matrix} \text{Wishart in} \\ \text{scalar case} \end{matrix}$$

Exercise: Let  $C \in \mathbb{R}^p$ ,  $A \sim \text{Wish}(\Sigma, m)$ ,  $\Sigma \in \mathbb{R}^{p \times p}$ .

$$C' A C \sim ?$$

$$C' A C \stackrel{?}{\sim} \text{Wish}(C' \Sigma C, m).$$

But now,  $C' A C > 0$  because  $\Sigma$  positive definite.

$$\Rightarrow \text{Wish}(C' \Sigma C, m) \sim (C' \Sigma C) \chi^2(m).$$

$$\Rightarrow C' A C \sim C' \Sigma C \chi^2(m).$$

Theorem:  $\underline{x}_1, \dots, \underline{x}_n$  iid  $\sim N_p(\mu, \Sigma)$ .

$$\Rightarrow \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' \sim \text{Wish}(\Sigma, n-1).$$

Corollary:  $S = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' \sim \text{Wish}\left(\frac{1}{n-1}\Sigma, n-1\right).$  distribution of  $S$

$$\hat{\Sigma} \sim \text{Wish}\left(\frac{1}{n}\Sigma, n-1\right).$$

We'll summarize all the informations about (the estimators) in a single theorem:

Theorem:  $\underline{x}_1, \dots, \underline{x}_n$  iid  $\sim N_p(\mu, \Sigma)$ . Then:

- 1)  $\bar{\underline{x}} \sim N_p(\mu, \frac{1}{n}\Sigma)$ ,
- 2)  $(n-1)S \sim \text{Wish}(\Sigma, n-1)$ ,
- 3)  $\bar{\underline{x}} \perp\!\!\!\perp S$  (stoch. independent).

distributions of  $\bar{\underline{x}}, S$   
in the multivariate Gaussian setting

Indeed, we have another theorem:

Theorem:  $\bar{\underline{x}}$  and  $S$  are sufficient statistics. [?]

▷ Useful transformations for making data more Gaussian:



If  $x$  is a proportion  $\in [0, 1]$ , a useful transformation is:

$$\log \frac{x}{1-x} =: \text{Logit}(x)$$

Asymptotic results for (Gaussian) r.v.s:

1) LLN: (Law of Large Numbers)

$\underline{x}_1, \dots, \underline{x}_n$ , random vectors iid s.t.  $E[\underline{x}_i] = \mu$ ,

$\text{cov}(\underline{x}_i) = \Sigma$  exist. Then:

$$\left[ \underline{x} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i \xrightarrow{P} \mu \right] \text{ as } n \rightarrow \infty.$$

$$\left[ S = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' \xrightarrow{P} \Sigma \right] \text{ as } n \rightarrow \infty.$$

## 2) CLT: (Teorema Centrale del Limite)

$\underline{x}_1, \dots, \underline{x}_n, \dots$  rv iid s.t.  $E[\underline{x}_i] = \mu$ ,  $\text{Cov}(\underline{x}_i) = \Sigma$  exist,  
then:  $\sqrt{n}(\bar{\underline{x}} - \mu) \sim N_p(\underline{0}, \Sigma)$ .  
 ↳ asymptotically normal

Meaning: for large  $n$ , one can approximate the distribution of  $\sqrt{n}(\bar{\underline{x}} - \mu)$  with a  $N_p(\underline{0}, \Sigma)$ .

In practice: for large  $n$ ,  $\bar{\underline{x}} \sim N_p(\mu, \frac{1}{n}\Sigma)$  approximately.

### Lab Session:

27/3/20

Lab 2: quick explanation.

↳ look at the code yourself.

### Lab 3: PCA.

### Inference for the mean $\mu$ :

30/3/20

We want to use the partial information we have in order to infer the data of the population.

$\mathbb{R}^p \ni \underline{x}_1, \dots, \underline{x}_n$  2. vectors iid, with  $E[\underline{x}_i] = \mu$  and  $\text{Cov}(\underline{x}_i) = \Sigma$ ,  $\det \Sigma > 0$ .

$$\bar{\underline{x}} = \frac{1}{n} \sum_{i=1}^n \quad (\text{point}) \text{ estimator for } \mu.$$

### I) Large $n$ : ( $n \gg p$ )

$$i) \sqrt{n}(\bar{\underline{x}} - \mu) \underset{\text{approx.}}{\sim} N_p(\underline{0}, \Sigma) \quad (\text{CLT})$$

$$\Rightarrow (\sqrt{n}(\bar{\underline{x}} - \mu))' \Sigma^{-1} (\sqrt{n}(\bar{\underline{x}} - \mu)) \sim \chi^2(p),$$

$$\text{or } n(\bar{\underline{x}} - \mu)' \Sigma^{-1} (\bar{\underline{x}} - \mu) \sim \chi^2(p).$$

The problem is that we don't know  $\Sigma$  as well.

If  $\Sigma$  is unknown, we have that:

$$ii) S \xrightarrow{P} \Sigma \quad \text{as } n \rightarrow \infty. \quad (\text{LLN})$$

Hence, for large  $n$ ,

$$n(\bar{X} - \mu)' S^{-1} (\bar{X} - \mu) \sim \chi^2(p).$$

distribution of the pivotal statistics with large  $n$

This is called a pivotal statistic: it's not a quantity that we know, but its distribution is known; we can make inference on it!

Consequences:

$$\Pr[n(\bar{X} - \mu)' S^{-1} (\bar{X} - \mu) \leq \chi_{\alpha}^2(p)] = 1 - \alpha \quad \text{if } \alpha \in (0, 1).$$

$$\Rightarrow \Pr[d_{S^{-1}}^2(\bar{X}, \mu) \leq \chi_{\alpha}^2(p)] = 1 - \alpha. \quad (*)$$

Define now the following sets:

$$\mathcal{E}_{S^{-1}}^{\alpha}(\mu) := \left\{ \underline{x} \in \mathbb{R}^p : d_{S^{-1}}^2(\bar{X}, \mu) \leq \chi_{\alpha}^2(p) \right\}$$

$$\mathcal{E}_{S^{-1}}^{\alpha}(\bar{X}) := \left\{ \eta \in \mathbb{R}^p : d_{S^{-1}}^2(\eta, \bar{X}) \leq \chi_{\alpha}^2(p) \right\}$$

Two ellipses, same axes, same radius. One centered on  $\mu$ , the other on  $\bar{X}$ . Their relationship is:

$$\bar{X} \in \mathcal{E}_{S^{-1}}^{\alpha}(\mu) \Leftrightarrow \mu \in \mathcal{E}_{S^{-1}}^{\alpha}(\bar{X}).$$

So  $(*)$  is equivalent to:

$$\Pr[\bar{X} \in \mathcal{E}_{S^{-1}}^{\alpha}(\mu)] = 1 - \alpha.$$

$\updownarrow$   
random fact

$$\Pr[\mu \in \mathcal{E}_{S^{-1}}^{\alpha}(\bar{X})] = 1 - \alpha.$$

$\hookrightarrow$   
random fact

Confidence region (level  $\alpha$ ) for  $\mu$ :

$$\begin{aligned} CR_{1-\alpha}(\mu) &= \left\{ \eta \in \mathbb{R}^p : \eta \in \mathcal{E}_{S^{-1}}^{\alpha}(\bar{X}) \right\} = \\ &= \left\{ \eta \in \mathbb{R}^p : n(\eta - \bar{X})' S^{-1}(\eta - \bar{X}) \leq \chi_{\alpha}^2(p) \right\}. \end{aligned}$$

Basics of inferential statistics:

Test:  $H_0: \mu = \mu_0$  vs.  $H_1: \mu \neq \mu_0$ .

$\alpha \in (0, 1)$  (small) : level of confidence of the test.

Consider the statistics:

$$T_0^2 := n (\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0).$$

If  $H_0$  is true, then  $T_0^2 \stackrel{(H_0)}{\sim} \chi^2(p)$ .

$\Rightarrow$  we have proof against  $H_0$  if  $\bar{X}$  is very far from  $\mu_0$ .

How big should that distance be?

Reject if  $T_0^2 > \chi_{\alpha}^2(p)$ . Either it's wrong the assumption, or we are observing a "miracle" (something that happens  $\alpha$  times every 1).

$$\boxed{\text{Rejection Region } \alpha = \{ T_0^2 > \chi_{\alpha}^2(p) \}}.$$

rejection region (level  $\alpha$ )  
for  $H_0$

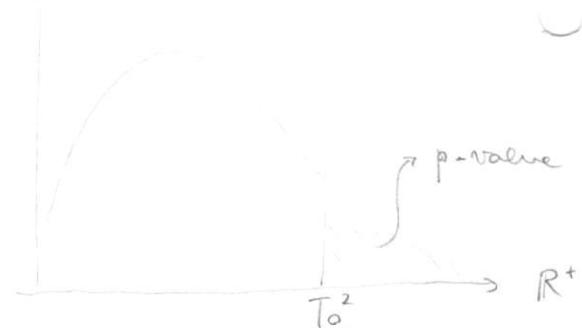
Normally we don't fix  $\alpha$ , but instead we look for the p-value of the test.

Definition:

p-value: area of  $\chi^2(p)$  at the right of  $T_0^2$ .

Testing with p-value: from the definition,

$$\boxed{\text{p-value} \leq \alpha \Leftrightarrow T_0^2 > \chi_{\alpha}^2(p)}.$$



• p-value small: we can reject  $H_0$ .

• p-value large: we cannot reject  $H_0$ .

History of 5% - 1%: values computed by Fisher on the tables.

$\hookrightarrow$  only a "historical reason". There is no absolute definition of what is large and small.

All this inference is based on knowing the distribution of the pivotal statistics. What if the sample is not large enough? We want to find another pivotal statistics for small  $n$ .

## II) Small $n$ :

Fisher's idea: convert the variables into Gaussians.

• 1905: "Gosset", brewer of the Guinness factory in Dublin, aka "Student"  $\rightarrow$  father of the "Student-t".

Key Assumption:  $\boxed{\underline{X}_1, \dots, \underline{X}_n \text{ iid } \sim N_p(\mu, \Sigma)}, \det \Sigma > 0$ .  
 $\hookrightarrow$  key fact

What is the distribution of

$n (\bar{X} - \mu)' S^{-1} (\bar{X} - \mu)$ ? We need to introduce Fisher's distribution.

Definition: Let  $Y \sim \chi^2(n)$ ,  $W \sim \chi^2(m)$ ,  $Y \perp\!\!\!\perp W$ .

$$\Rightarrow \frac{Y/n}{W/m} \sim F(n, m).$$

Fisher's Distribution

↳ coming from Fisher

Observations: properties of Fisher's distributions.

$$1) t \sim t(n) \Rightarrow t = \frac{Z}{\sqrt{W/n}}, \begin{cases} Z \sim N(0, 1) \\ W \sim \chi^2(n) \\ Z \perp\!\!\!\perp W \end{cases}$$

Taking:

$$t^2 = \frac{Z^2}{W/n}, \begin{cases} Z^2 \sim \chi^2(1) \\ W \sim \chi^2(n) \\ Z^2 \perp\!\!\!\perp W \end{cases} \Rightarrow t^2 \sim F(1, m). \quad \begin{matrix} \text{Fisher with} \\ n=1 \text{ d.o.f.} \end{matrix}$$

$$2) F(n, m) \xrightarrow{m \rightarrow \infty} \frac{1}{n} \chi^2(n) \quad (\text{convergence in distribution}) \quad \begin{matrix} \text{Fisher with} \\ m \rightarrow \infty \text{ d.o.f's} \end{matrix}$$

Indeed,  $\frac{Y/n}{W/m}, Y \sim \chi^2(n), W \sim \chi^2(m), Y \perp\!\!\!\perp W$ .

$$\Rightarrow W = \sum_{i=1}^m z_i^2, z_1, \dots, z_m \text{ iid } \sim N(0, 1).$$

$$E[z_i^2] = 1 \Rightarrow \frac{1}{m} \sum_{i=1}^m z_i^2 \rightarrow 1 \text{ as } m \rightarrow \infty \quad (\text{LLN})$$

$$\Rightarrow \frac{W}{m} \rightarrow 1 \text{ as } m \rightarrow \infty \Rightarrow \frac{Y/n}{W/m} \rightarrow \frac{1}{n} Y \sim \frac{1}{n} \chi^2(n) \text{ as } m \rightarrow \infty.$$

Exercise: use R to make plots of different F distributions.

Theorem: Hotelling, 1931.

Suppose  $\underline{X} \sim N_p(\mu, \Sigma)$ ,  $\det \Sigma > 0$ ,  $W \sim \text{Wish}(\frac{1}{m} \Sigma, m)$ ,  $\Sigma$  p×p,  
 $\underline{X} \perp\!\!\!\perp W$ . Then:

$$\frac{m-p+1}{mp} (\underline{X} - \mu)' W^{-1} (\underline{X} - \mu) \sim F(p, m-p+1). \quad \begin{matrix} \text{Hotelling's} \\ \text{Theorem} \end{matrix}$$

Corollary: let  $X_1, \dots, X_n$  iid  $\sim N_p(\mu, \Sigma)$ .

$$\Rightarrow m (\bar{\underline{X}} - \mu)' S^{-1} (\bar{\underline{X}} - \mu) \sim \frac{(n-1)p}{n-p} F(p, n-p). \quad \begin{matrix} \text{distribution of} \\ \text{fimital statistics} \\ \text{with Gaussianity} \end{matrix}$$

Proof (corollary):

$$\sqrt{n}(\bar{X} - \mu) \sim N_p(0, \Sigma), \quad S \sim \text{Wish}(\frac{1}{n-1}\Sigma, n-1), \quad \bar{X} \perp\!\!\!\perp S.$$

Use Hotelling's thm. ( $m = n-1$ ):

$$\frac{n-1+p+1}{(n-1)p} n (\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu) \sim F(p, n-1-p+1)$$

$$\Rightarrow \frac{n-p}{(n-1)p} n (\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu) \sim F(p, n-p)$$

$$\Leftrightarrow n (\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu) \sim \frac{(n-1)p}{n-p} F(p, n-p)$$

(end)

We have discovered the distribution of  $T^2 := n (\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu)$ !  
This is the pivotal statistics for making inference with  $\mu$ .

• Confidence Region:

$$\Pr[n (\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu)] \leq \frac{(n-1)p}{n-p} F_\alpha(p, n-p) = 1-\alpha$$

$d_{S^{-1}}^2(\bar{X}, \mu)$

$$\Rightarrow \Pr[\mu \in \mathcal{E}_{S^{-1}}^\alpha(\bar{X})] = 1-\alpha, \text{ with } \alpha \in (0, 1) \text{ and}$$

$$\mathcal{E}_{S^{-1}}^\alpha(\bar{X}) = \left\{ \eta \in \mathbb{R}^p : d_{S^{-1}}^2(\bar{X}, \eta) \leq \frac{(n-1)p}{n-p} F_\alpha(p, n-p) \right\}.$$

$$\Rightarrow \boxed{\text{CR}_{1-\alpha}(\mu) = \mathcal{E}_{S^{-1}}^\alpha(\bar{X})}.$$

same Mahalanobis distance as before, but different radius!

Observation:

$$\underbrace{\frac{(n-1)p}{n-p} F_\alpha(p, n-p)}_{\downarrow 1} \rightarrow \chi_\alpha(p)$$

as  $n \rightarrow \infty$ .

convergence of the two pivotal quantities as  $n \rightarrow \infty$

• Testing:  $H_0: \mu = \mu_0$  vs.  $H_1: \mu \neq \mu_0$  at level  $\alpha \in (0, 1)$ .

$$T_0^2 = n (\bar{X} - \mu_0)^T S^{-1} (\bar{X} - \mu_0).$$

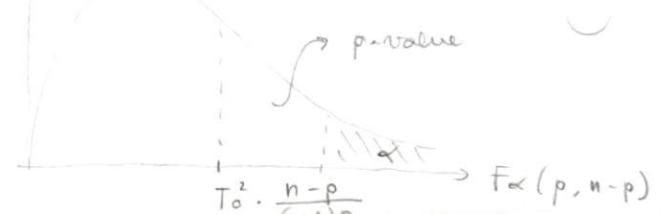
$$\text{If } H_0 \text{ is true} \Rightarrow T_0^2 \stackrel{(H_0)}{\sim} \frac{(n-1)p}{n-p} F(p, n-p).$$

$$\text{Reject if } T_0^2 > \frac{(n-1)p}{n-p} F_\alpha(p, n-p).$$

Remark: this time, the p-value is the area at the right of:

$$T_0^2 \cdot \frac{n-p}{(n-1)p}.$$

Don't forget the rescaling factor!



Example:  $\underline{X}_1, \dots, \underline{X}_{10}$  iid  $\sim N_2(\mu, \Sigma)$ .

Suppose that  $\bar{\underline{X}} = \underline{0}$ , and  $S = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ .

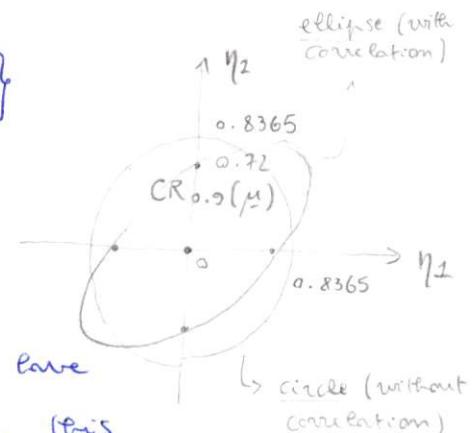
$$CR_{1-\alpha}(\mu) = \{\eta \in \mathbb{R}^2 : 10 \eta' I \eta \leq \frac{9.2}{8} F_{\alpha}(2, 8)\}.$$

Choose  $\alpha = 10\%$ .  $\Rightarrow \{1-\alpha = 0.9\}$   
 $F_{0.1}(2, 8) = 3.11$ .

$$\Rightarrow CR_{0.9}(\mu) = \{\eta \in \mathbb{R}^2 : 10(\eta_1^2 + \eta_2^2) \leq 6.997\}$$

$$= \{\eta \in \mathbb{R}^2 : \eta_1^2 + \eta_2^2 \leq 0.6997\}.$$

Observe that  $\mu$  is a parameter, and what is random is the circle we have drawn. We have tossed the coin and today we have observed this ellipse. This ellipse is right 90% of the times.



Assume now  $\bar{\underline{X}} = \underline{0}$ ,  $S = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$ .

$$\Rightarrow CR_{0.9}(\mu) = \{\eta \in \mathbb{R}^2 : 10 \eta' S^{-1} \eta \leq \frac{9.2}{8} \cdot 3.11\}$$

$$S^{-1} = \begin{bmatrix} 4/3 & -2/3 \\ -2/3 & 4/3 \end{bmatrix} \Rightarrow 10 \eta' S^{-1} \eta = \frac{10 \cdot 4}{3} (\eta_1^2 - \eta_1 \eta_2 + \eta_2^2).$$

$$\Rightarrow CR_{0.9}(\mu) = \{\eta \in \mathbb{R}^2 : \eta_1^2 - \eta_1 \eta_2 + \eta_2^2 \leq \frac{3}{4} \cdot 0.6997\}.$$

↳ we have an ellipse this time! It changed shape because  $\Sigma$  changed. The more the two components are correlated, the more the two variables will be squished to a bisector.

#### univariate vs. multivariate CIs:

Now,  $\bar{\underline{X}} = \begin{pmatrix} \bar{X}_1 \\ \bar{X}_2 \end{pmatrix} = \underline{0}$ ,  $S = \begin{pmatrix} S_{11} = 1 & S_{12} \\ S_{21} & S_{22} = 1 \end{pmatrix}$ .

Consider the univariate confidence intervals:

$$\text{CI}_{0.9}(\mu_1) = [\bar{X}_1 \pm t_{0.95}(9) \sqrt{\frac{1}{10}}] = [\pm 0.5796].$$

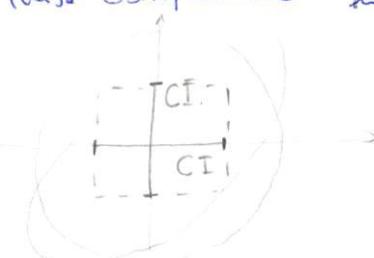
$$\text{CI}_{0.9}(\mu_2) = [\pm 0.5796]$$

Here we didn't use the covariance for these Confidence Intervals!

CI's are independent from  $S$ .

They seem better intervals than before!

But here we have something fishy:



$$\Pr[\mu_1 \in CI_{\alpha/2}(\bar{X}_1), \mu_2 \in CI_{\alpha/2}(\bar{X}_2)] = 0.81$$

$$\Pr[\bar{X}_1 \leq \bar{X}_2] = \Pr[\mu_1 \in CI_{\alpha/2}(\bar{X}_1)] \Pr[\mu_2 \in CI_{\alpha/2}(\bar{X}_2)] = (0.9)^2 < 0.9.$$

Hence, the probability for  $(\mu_1, \mu_2)$  to fall inside the cartesian square is actually less than 0.9! In other words, we are giving a multi-variate CI with a level of  $(1-\alpha)^2$ , smaller: the true  $\alpha$  here is  $\alpha = 19\%$ .

Quick presentation: 7th of April.

31/3/20

First slide: title of the dataset.

Second slide: faces.

Description of the dataset, what we have already done.

Suggestion of the professor: follow the face of the professor's notes, but study on the book. Don't study just the notes!

### Confidence Intervals for linear combinations of the mean:

We are still in the small  $n$  sample:

$$X_1, \dots, X_n \text{ iid } \sim N_p(\mu, \Sigma)$$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \text{ estimator for } \mu.$$

I) One-at-the-time CI for  $\mu$ :

let  $a \in \mathbb{R}$ . What is  $[CI_{1-\alpha}(a' \mu)]$ ?

Estimator for  $a' \mu$ :  $a' \bar{X}$ .

$$a' \bar{X} \sim N_1(a' \mu, \frac{1}{n} a' \Sigma a), \text{ i.e. } \frac{a' \bar{X} - a' \mu}{\sqrt{a' \Sigma a}} \sqrt{n} \sim N_1(0, 1).$$

We don't know  $\Sigma$ , but we can estimate it with  $S$ :

$$(n-1)S \sim \text{Wishart}(\Sigma, n-1)$$

$$\Rightarrow (n-1)a' S a \sim (a' \Sigma a) \chi^2(n-1).$$

Moreover,  $\bar{X} \perp\!\!\!\perp S$ . Taking a ratio between a Gaussian and a normalized squared Wishart ( $\chi^2$ ), we get a Student-t:  
(by definition of Student-t)

$$\frac{a' \bar{X} - a' \mu}{\sqrt{a' \Sigma a}} \sqrt{n} \sim N(0, 1)$$

$$\frac{(n-1)a' S a}{\sqrt{(n-1)(a' \Sigma a)}} \sim \sqrt{\frac{\chi^2(n-1)}{n-2}}$$

Simplifying:

$$\frac{a' \bar{X} - a' \mu}{\sqrt{a' S a}} \sqrt{n} \sim t(n-1).$$

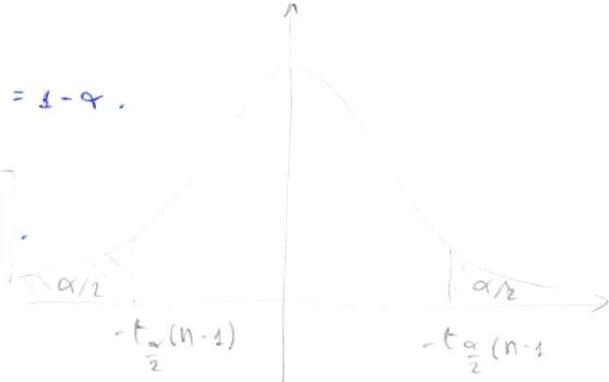
pivotal quantity for one-at-the-time CI of  $a' \mu$

We have obtained our pivotal quantity! For example,

$$\Pr \left[ \frac{|a' \bar{X} - a' \mu|}{\sqrt{a' S a}} \sqrt{n} < t_{\frac{\alpha}{2}}(n-1) \right] = 1-\alpha, \quad \alpha \in (0, 1).$$

$$\Rightarrow \Pr [a' \mu \in \left[ a' \bar{X} \pm t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{a' S a}{n}} \right]] = 1-\alpha.$$

$$\Rightarrow CI_{1-\alpha}(a' \mu) = \left[ a' \bar{X} \pm t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{a' S a}{n}} \right].$$



Example:  $\underline{a} = (0 \dots 0 \underset{i}{1} 0 \dots 0 \dots 0)^T, i=1 \dots p.$

$$\Rightarrow \underline{a}' \mu = \mu_i.$$

$$CI_{1-\alpha}(\mu_i) = \left[ \bar{X}_i \pm t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{s_{ii}}{n}} \right].$$

If  $\underline{a} = (0 \dots 0 \underset{i}{1} 0 \dots 0 \underset{j}{1} 0 \dots 0)^T \Rightarrow \underline{a}' \mu = \mu_i - \mu_j$ .

$$\Rightarrow CI_{1-\alpha}(\mu_i - \mu_j) = \left[ \bar{X}_i - \bar{X}_j \pm t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{s_{ii} - 2s_{ij} + s_{jj}}{n}} \right]. \quad \square$$

We proved that:

$$\Pr [a' \mu \in \left[ a' \bar{X} \pm t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{a' S a}{n}} \right]] = 1-\alpha, \quad \forall a \in \mathbb{R}^p.$$

one-at-the-time  
CI for  $a' \mu$

Aside: testing  $H_0: a' \mu = \delta_0$ .

$$1) H_1: a' \mu \neq \delta_0. \quad \text{If } H_0 \text{ is true, } \frac{a' \bar{X} - \delta_0}{\sqrt{a' S a}} \sqrt{n} \sim t(n-1).$$

So we reject if

$$\frac{a' \bar{X} - \delta_0}{\sqrt{a' S a}} \sqrt{n} > t_{\alpha/2}(n-1).$$

2) If  $H_1$  is  $a' \mu \neq \delta_0$ , we reject at level  $\alpha \in (0, 1)$  if:

$$\frac{|a' \bar{X} - \delta_0|}{\sqrt{a' S a}} \sqrt{n} > t_{\alpha/2}(n-1).$$



We are not learning these formulas by heart, we are obtaining them from pivotal distributions!

Question: did we prove that:

$$(*) \Pr \left[ \underline{a}' \mu \in \left[ \underline{a}' \bar{x} + t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{\underline{a}' S \underline{a}}{n}} \right], \forall \underline{a} \in \mathbb{R}^p \right] = 1-\alpha ?$$

[No, this is totally wrong. That's why we are computing "one CI at-the-time".]

What is wrong is the term  $t_{\frac{\alpha}{2}}(n-1)$ : we need to construct a different pivotal quantity for this new problem.

II). Simultaneous CI for linear combinations for  $\mu$ :

We want something like this:

$$\Pr \left[ \underline{a}' \mu \in \left[ \underline{a}' \bar{x} + (?) \sqrt{\frac{\underline{a}' S \underline{a}}{n}} \right], \forall \underline{a} \in \mathbb{R}^p \right] = 1-\alpha.$$

We are interested in simultaneous inference, that doesn't have to be mixed with one-at-the-time interval.

We need a small excusus on linear algebra.

Example: coin toss.

$$X_1 = \begin{cases} 0, & 1/2 \\ 1, & 1/2 \end{cases}, \quad X_2 = \begin{cases} 0, & 1/2 \\ 1, & 1/2 \end{cases} \quad X_1, \dots, X_n \text{ iid } \text{Bern}\left(\frac{1}{2}\right).$$

$$1) \Pr[X_i = 1] = \frac{1}{2} \quad \forall i = 1, \dots, 100. \quad \text{one-at-the-time}$$

$$2) \Pr[X_i = 1, \forall i = 1, \dots, 100] \neq \frac{1}{2}. \quad \text{simultaneously}$$

$$\hookrightarrow = \left(\frac{1}{2}\right)^{100}.$$

A little excusus in Algebra:

$$\text{Let } \underline{b}, \underline{d} \in \mathbb{R}^p. \quad \frac{\langle \underline{b}, \underline{d} \rangle}{\|\underline{b}\| \|\underline{d}\|} = \cos(\vartheta).$$

$$\Rightarrow \frac{(\langle \underline{b}, \underline{d} \rangle)^2}{\|\underline{b}\|^2 \|\underline{d}\|^2} \leq 1. \quad \text{Equality holds if } \underline{b} \in \text{span}(\underline{d}), \text{ i.e. } \underline{b} \text{ or } \underline{d}.$$

$$\text{Hence, } (\underline{b}' \underline{d})^2 \leq \left( \sum_i b_i^2 \right) \left( \sum_i d_i^2 \right) \quad \text{and equality holds if } \underline{b} \in \text{span}(\underline{d}).$$

Writing it in a different way,

$$\int fg \leq (\int f^2)(\int g^2). \quad \text{Cauchy-Schwarz inequality}$$

[all of the statistics in  $\mathbb{R}^p$  can be transfert in a Hilbert space]

Take now  $B$   $p \times p$  positive definite matrix.

We have that:

$$(\underline{b}' \underline{d})^2 \leq (\underline{b}' B \underline{b}) (\underline{d}' B^{-1} \underline{d})$$

extended Cauchy-Schwarz  
(CS)

Proof:

$$\begin{aligned} \underline{b}' \underline{d} &= \underline{b}' B^{1/2} B^{-1/2} \underline{d} \\ \Rightarrow (\underline{b}' \underline{d})^2 &= (\underline{b}' B^{1/2} B^{-1/2} \underline{d})^2 \stackrel{\text{C.S.}}{\leq} (\underline{b}' B^{1/2} B^{1/2} \underline{b}) (\underline{d}' B^{-1/2} B^{-1/2} \underline{d}') \\ &= (\underline{b}' B \underline{b}) (\underline{d}' B \underline{d}). \end{aligned}$$

Equality holds if  $\underline{b}' B^{1/2} \in \text{span}(B^{-1/2} \underline{d})$ , or  $\underline{b} \in \text{span}(B^{-1} \underline{d})$ .

("maximum lemma")

(end)

Proposition: Let  $B$   $p \times p$  pos. def.,  $\underline{d} \in \mathbb{R}^p$ . Then:

$$\max_{\underline{x} \in \mathbb{R}^p, \underline{x} \neq \underline{0}} \frac{(\underline{x}' \underline{d})^2}{\underline{x}' B \underline{x}} = \underline{d}' B^{-1} \underline{d}. \quad \text{maximum lemma}$$

Proof:  $(\underline{x}' \underline{d})^2 \leq (\underline{x}' B \underline{x})(\underline{d}' B^{-1} \underline{d})$ .

If  $\underline{x} \neq \underline{0} \Rightarrow \underline{x}' B \underline{x} > 0$ .

$$\Rightarrow \frac{(\underline{x}' \underline{d})^2}{\underline{x}' B \underline{x}} \leq \underline{d}' B^{-1} \underline{d} \quad \forall \underline{x} \neq \underline{0}.$$

Equality holds if  $\underline{x} \in \text{span}(B^{-1} \underline{d})$ .

(end)

Back to statistics:

$$\underline{x}_1, \dots, \underline{x}_n \text{ iid } \sim N_p(\mu, \Sigma), \quad \bar{\underline{x}} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i.$$

Let  $\underline{\alpha} \in \mathbb{R}^p$ , take  $\frac{(\underline{\alpha}' (\bar{\underline{x}} - \mu))^2}{\underline{\alpha}' \underline{\Sigma} \underline{\alpha}}$ . By maximum lemma,

$$\max_{\underline{\alpha} \in \mathbb{R}^p} n \frac{(\underline{\alpha}' (\bar{\underline{x}} - \mu))^2}{\underline{\alpha}' \underline{\Sigma} \underline{\alpha}} = n \underbrace{(\bar{\underline{x}} - \mu)' \underline{\Sigma}^{-1} (\bar{\underline{x}} - \mu)}_{\parallel} \sim \frac{(n-1)p}{n-p} F(p, n-p).$$

$\Rightarrow$  We have our pivotal statistics! Notice that it's the same elliptical statistics that we use to build the confidence Region. Therefore:

$$\Pr \left[ \frac{|\underline{\alpha}' (\bar{\underline{x}} - \mu)|}{\sqrt{\underline{\alpha}' \underline{\Sigma} \underline{\alpha}}} \sqrt{n} \leq c, \quad \forall \underline{\alpha} \in \mathbb{R}^p, \underline{\alpha} \neq \underline{0} \right] = 1-\alpha.$$

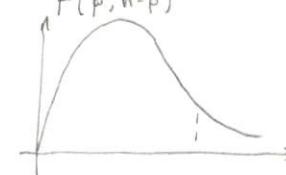
II  $\rightarrow$  key equality  $(*)$

$$\Pr \left[ \max_{\substack{\underline{\alpha} \in \mathbb{R}^p \\ \underline{\alpha} \neq \underline{0}}} \frac{|\underline{\alpha}' (\bar{\underline{x}} - \mu)|}{\sqrt{\underline{\alpha}' \underline{\Sigma} \underline{\alpha}}} \sqrt{n} \leq c \right] = 1-\alpha.$$

(\*) if the quantity is smaller  $\forall \underline{\alpha}$ , then it will be smaller also for the  $\underline{\alpha} \in \mathbb{R}^p$  generating the biggest quantity.

Taking the square (monotone function):

$$\Pr \left[ \max_{\substack{\alpha \in \mathbb{R}^p \\ \alpha \neq 0}} \frac{(\alpha' (\bar{X} - \mu))^2}{\alpha' \Sigma \alpha} n \leq c^2 \right] = 1 - \alpha$$



$$\sim \frac{(n-1)p}{n-p} F(p, n-p) \Rightarrow c^2 = \frac{(n-1)p}{n-p} F_\alpha(p, n-p).$$

We can conclude that:

$$\Pr \left[ \frac{|\alpha' (\bar{X} - \mu)|}{\sqrt{\alpha' \Sigma \alpha}} \sqrt{n} \leq \sqrt{\frac{(n-1)p}{n-p} F_\alpha(p, n-p)} \quad \forall \alpha \in \mathbb{R}^p \right] = 1 - \alpha.$$

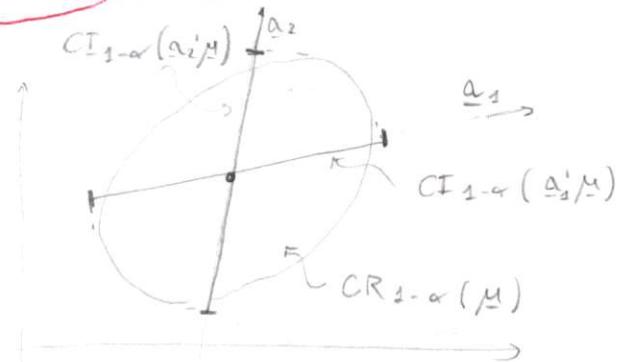
The quantile is a Fisher, not a Student-t!

$$CI_{1-\alpha}(\alpha' \mu) = \left[ \alpha' \bar{X} \pm \sqrt{\frac{(n-1)p}{n-p} F_\alpha(p, n-p)} \sqrt{\frac{\alpha' \Sigma \alpha}{n}} \right].$$

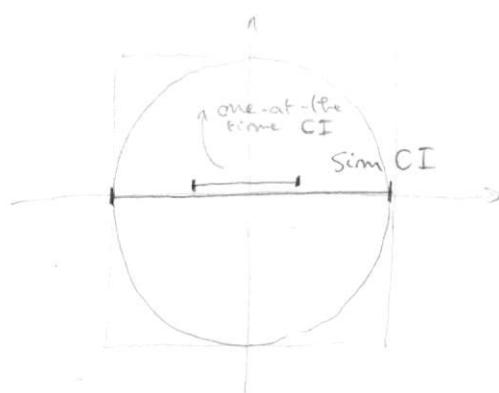
since it's a vector  
CI for  $\alpha' \mu$

### Geometry of the CI:

The simultaneous CI is the projection of the Confidence Region of  $\mu$  along all the directions  $\alpha$ .



Comment: today, with Big Data, we want to estimate a large number of parameters. Simultaneous CI is hence very used.



We have two "extremes":

- simultaneous cover all possible linear combinations: longer CIs.
- one-at-the-time CIs: too specific.

What if we specify just a finite number of them? Your client wants to see a smaller CI, so we are tempted to sell them the one-at-the-time CI. The intermediate method is given by Bonferroni.